

434 Appendix

435 A Theoretical Proof

436 We provide a detailed derivation and proof of the least-squares pose regression process in the text.
 437 The camera position parameters θ and t can be obtained by solving the following loss function.

$$\xi = \arg \min_{\theta, t} \sum_{i=1}^n S_i(\mathbf{R}\mathbf{p}'_i + \mathbf{t} - \hat{\mathbf{p}}_i)^2, \mathbf{R} = \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix} \quad (8)$$

438 The relationship between the corresponding points \mathbf{p}'_i and $\hat{\mathbf{p}}_i$ in point sets \mathbf{p}' and $\hat{\mathbf{p}}$ can be expressed
 439 as $\mathbf{p}'_i = \mathbf{R}\hat{\mathbf{p}}_i + \mathbf{t}$. To simplify the problem, we calculate the centroids \mathbf{g}' and $\hat{\mathbf{g}}$ of point sets \mathbf{p}' and
 440 $\hat{\mathbf{p}}$, respectively. We then subtract the centroids from each point in their respective sets, resulting in
 441 $\mathbf{q}'_i = \mathbf{p}'_i - \mathbf{g}'$ and $\hat{\mathbf{q}}_i = \hat{\mathbf{p}}_i - \hat{\mathbf{g}}$. This simplifies the problem to:

$$\begin{aligned} \xi &= \arg \min_{\theta} \sum_{i=1}^n S_i(\mathbf{R}\mathbf{q}'_i - \hat{\mathbf{q}}_i)^2 \\ &= \arg \min_{\theta} \sum_{i=1}^n S_i(\mathbf{R}\mathbf{q}'_i - \hat{\mathbf{q}}_i)^t (\mathbf{R}\mathbf{q}'_i - \hat{\mathbf{q}}_i) \\ &= \arg \min_{\theta} \sum_{i=1}^n S_i(\mathbf{q}'_i{}^t \mathbf{q}'_i) + \hat{\mathbf{q}}_i{}^t \hat{\mathbf{q}}_i - 2\hat{\mathbf{q}}_i{}^t \mathbf{R}\mathbf{q}'_i \end{aligned} \quad (9)$$

442 Therefore, minimizing ξ is equivalent to maximizing $\sum_{i=1}^n S_i(\hat{\mathbf{q}}_i{}^t \mathbf{R}\mathbf{q}'_i)$. Assuming that $\mathbf{H} =$
 443 $\sum_{i=1}^n S_i(\mathbf{q}'_i \hat{\mathbf{q}}_i{}^t)$, then

$$\sum_{i=1}^n S_i(\hat{\mathbf{q}}_i{}^t \mathbf{R}\mathbf{q}'_i) = \text{tr}(\sum_{i=1}^n \mathbf{R} S_i \mathbf{q}'_i \hat{\mathbf{q}}_i{}^t) = \text{tr}(\mathbf{R}\mathbf{H}) \quad (10)$$

444 Upon performing an SVD decomposition on \mathbf{H} , the result is $\mathbf{H} = \mathbf{U}\mathbf{\Lambda}\mathbf{V}^t$, where \mathbf{U} and \mathbf{V} are
 445 orthogonal matrices of dimensions 3×3 , and $\mathbf{\Lambda}$ is 3×3 diagonal matrix with non-negative elements.

446 Then we introduce an orthogonal matrix $\mathbf{X} = \mathbf{U}\mathbf{V}^t$, and observe that $\mathbf{X}\mathbf{H} = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^t$ is a symmet-
 447 ric positive-definite matrix. According to the properties of positive definite matrices, it holds that
 448 $\text{tr}(\mathbf{X}\mathbf{H}) \geq \text{tr}(\mathbf{B}\mathbf{X}\mathbf{H})$ for any 3×3 orthogonal matrix \mathbf{B} . Therefore, we obtain the maximum
 449 value of $\sum_{i=1}^n S_i(\hat{\mathbf{q}}_i{}^t \mathbf{R}\mathbf{q}'_i)$ when $\mathbf{R} = \mathbf{X}$, which implies that

$$\mathbf{R} = \mathbf{U}\mathbf{V}^t \quad (11)$$

450 The dense matching relationship between the graphs gives rise to the uniqueness of solution for \mathbf{R} .
 451 Subsequently, the relative displacement can be determined as

$$\mathbf{t} = \mathbf{g}' - \mathbf{R}\hat{\mathbf{g}} \quad (12)$$

452 By utilizing equations 11 and 12, we can obtain the camera location and azimuth through a straight-
 453 forward differentiable process, using the least-squares regression method.

454 B Localization qualitative results

455 In Figures 5 and 6, we present additional qualitative cross-view localization results on the KITTI,
 456 Ford multi-AV, VIGOR, and Oxford RobotCar datasets.

457 As pointed out in the main paper Section 4.6 of the main text, our method tends to exhibit poorer
 458 performance in terms of longitudinal localization accuracy. We present some examples of localiza-
 459 tion errors in Figure 5. In the visualization results it can be noticed that our algorithm demonstrates
 460 higher accuracy in lateral localization along the direction of travel, while performing poorly in lon-
 461 gitudinal localization, particularly in the Ford multi-AV dataset of suburban scenes. The displayed
 462 examples of localization failures reveal that the limited horizontal field of view and the absence of
 463 reference objects along the sides of the lane contribute to the localization errors. In such scenarios,
 464 our algorithm is still able to rely on lane features to accomplish lateral localization compared to LM.
 465 However, neither method performs well in longitudinal localization.

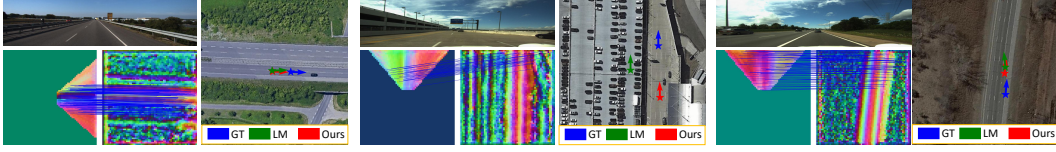


Figure 5: Failure cases. The left image sourced from KITTI dataset and the two on the right from the Ford multi-AV dataset. For each scene, the up left is the ground image, the bottom left denotes matching inliers and the right shows the satellite image and localization results. For visual simplicity, dense matching displays only partial inliers.

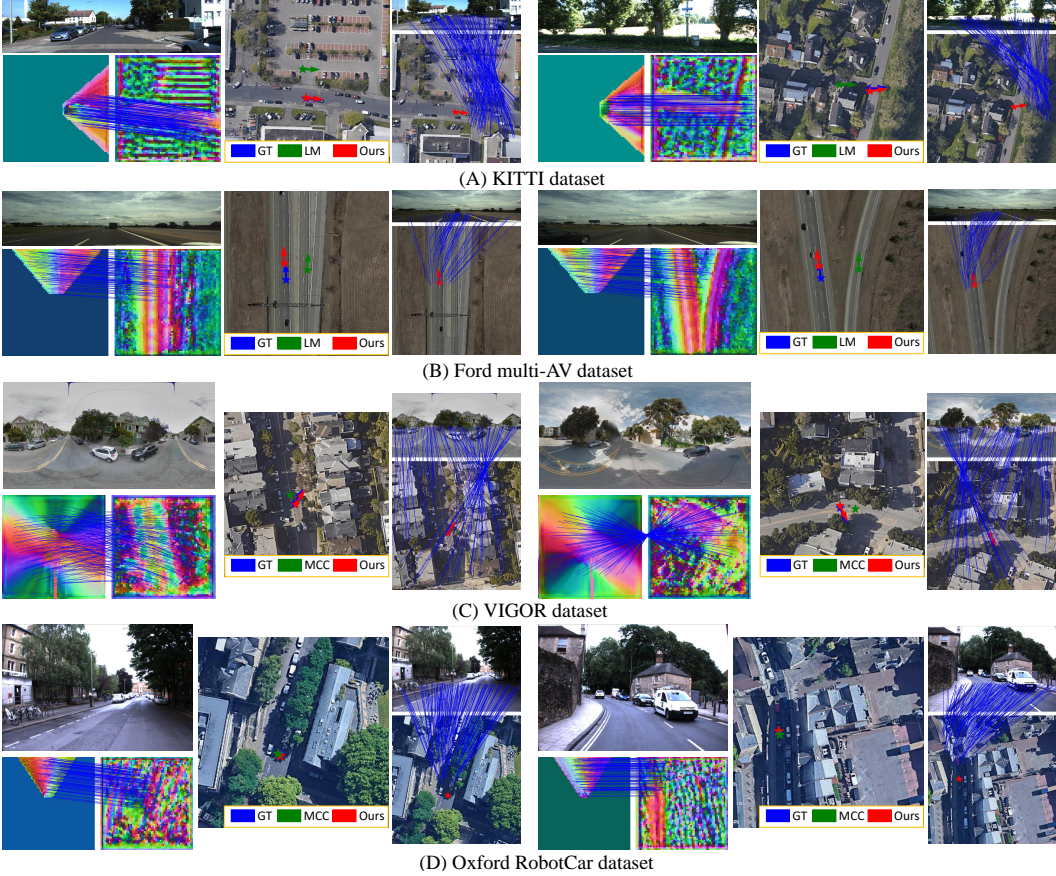


Figure 6: Qualitative results on the four datasets. For each scene, the up left is the ground image, the bottom left denotes matching inliers, the center shows the satellite image and localization results, and the right displays a visualization of the matched points mapped back to the original image. The central axis of a panoramic image represents its orientation. For visual simplicity, dense matching displays only partial inliers.

466 In addition, we have provided more visualizations for all four datasets. As discussed in Section 4 of
 467 the main text, our approach outperforms previous methods. It is important to note that, during testing
 468 on the VIGOR dataset, the MCC method cannot handle angle noise. Therefore, the visualization for
 469 MCC is based on a scenario where the angle noise is zero, while our visualization results are obtained
 470 in a scenario where the angle is completely unknown. Furthermore, to demonstrate the effectiveness
 471 of dense matching, we visualize the dense matching points mapped onto the original images in
 472 Figure 5. It can be observed that our pipeline establishes correspondences between satellite and
 473 ground-level image points, achieving precise localization of ground cameras in satellite images. This
 474 further confirms the efficacy of the architecture proposed in our work for cross-view localization
 475 using dense pixel flow fields.